

# A nonlinear aggregation type classifier

Alejandro Cholaquidis, Ricardo Fraiman, Juan Kalemkerian  
and Pamela Llop

September 10, 2015

## Abstract

We introduce a nonlinear aggregation type classifier for functional data defined on a separable and complete metric space. The new rule is built up from a collection of  $M$  arbitrary training classifiers. If the classifiers are consistent, then so is the aggregation rule. Moreover, asymptotically the aggregation rule behaves as well as the best of the  $M$  classifiers. The results of a small simulation are reported both, for high dimensional and functional data, and a real data example is analyzed.

*Keywords:* Functional data; supervised classification; non-linear aggregation.

## 1 Introduction

Supervised classification is still one of the hot topics for high dimensional and functional data due to the importance of their applications and the intrinsic difficulty in a general setup. In this context, there is a vast literature on classification methods which include: linear classification,  $k$ -nearest neighbors and kernel rules, classification based on partial least squares, reproducing kernels or depth measures. Complete surveys of the literature are the works by Baíllo et al. [1], Cuevas [13] and Delaigle and Hall [16]. In the book *Contributions in infinite-dimensional statistics and related topics* [7], there are also several recent advances in supervised and unsupervised classification. See for instance, Chapters 2, 5, 22 or 48, or directly, Chapter 1 of this issue (Bongiorno et al. [6]). In this context, very recently there have been of great interest to develop aggregation methods. In particular, there is a large list of linear aggregation methods like boosting (Breiman [8], Breiman [9]), random forest (Breiman [10], Biau et al. [3], Biau [5]), among others. All these methods exhibit an important improvement when combining a subset of classifiers to produce a new one. Most of the contributions to the aggregation literature have been proposed for nonparametric regression, a problem closely related to classification rules, which can be obtained just by plugging in the estimate of the regression function into the Bayes rule (see for instance, Yang [19] and Bunea et al. [11]). Model selection (select the optimal single model from a list of models), convex aggregation (search for the optimal convex combination of a given set of estimators),

and linear aggregation (select the optimal linear combination of estimators) are important contributions among a large list.

In the finite dimensional setup, Mojirsheibani [17] and [18] introduced a combined classifier showing strong consistency under somewhat hard to verify assumptions involving the Vapnik Chervonenkis dimension of the random partitions of the set of classifiers, which are non-valid in the functional setup. Very recently Biau et al. [4] introduced a new nonlinear aggregation strategy for the regression problem called COBRA, extending the ideas in Mojirsheibani [17] to the more general setup of nonparametric regression in  $\mathbb{R}^d$ . In the same direction but for the classification problem in the infinite dimensional setup, we extend the ideas in Mojirsheibani [17] to construct a classification rule which combines, in a nonlinear way, several classifiers to construct an optimal one. We point out that our rule allows to combine methods of very different nature, taking advantage of the abilities of each expert and allowing to adapt the method to different class of datasets. Even though our classifier allows aggregate experts of the same nature, the possibility of combine classifiers of different character, improves the use of existing rules as the bagged nearest neighbors classifier (see for instance Hall and Samworth [15]). As in Biau et al. [4], we also introduce a more flexible form of the rule which discards a small percentage  $\alpha$  of those preliminary experts that behaves differently from the rest. Under very mild assumptions, we prove consistency, obtain rates of convergence and show some optimality properties of the aggregated rule. To build up this classifier, we use the inverse function (see also Fraiman et al. [14]) of each preliminary experts which makes the proposal particularly well designed for high dimensional data avoiding the curse of dimensionality. It also performs well in functional data settings.

In Section 2 we introduce the new classifier in the general context of a separable and complete metric space which combines, in a nonlinear way, the decision of  $M$  experts (classifiers). A more flexible rule is also considered. In Section 3 we state our two main results regarding consistency, rates of convergence and asymptotic optimality of the classifier. Asymptotically, the new rule performs as the best of the  $M$  classifiers used to build it up. Section 4 is devoted to show through some simulations the performance of the new classifier in high dimensional and functional data for moderate sample sizes. A real data example is also considered. All proofs are given in the Appendix.

## 2 The setup

Throughout the manuscript  $\mathcal{F}$  will denote a separable and complete metric space,  $(X, Y)$  a random pair taking values in  $\mathcal{F} \times \{0, 1\}$  and  $\mu$  the probability measure of  $X$ . The elements of the training sample  $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ , are iid random elements with the same distribution as the pair  $(X, Y)$ . The regression function is denoted by  $\eta(x) = \mathbb{E}(Y|X = x) = \mathbb{P}(Y = 1|X = x)$ , the Bayes rule by  $g^*(x) = \mathbb{I}_{\{\eta(x) > 1/2\}}$  and the optimal Bayes risk by  $L^* = \mathbb{P}(g^*(X) \neq Y)$ .

In order to define our classifier, we split the sample  $\mathcal{D}_n$  into two subsamples  $\mathcal{D}_k = \{(X_1, Y_1), \dots, (X_k, Y_k)\}$  and  $\mathcal{E}_l = \{(X_{k+1}, Y_{k+1}), \dots, (X_n, Y_n)\}$  with  $l = n - k \geq 1$ . With  $\mathcal{D}_k$  we build up  $M$  classifiers  $g_{mk} : \mathcal{F} \rightarrow \{0, 1\}$ ,  $m = 1, \dots, M$  which we place in the vector  $\mathbf{g}_k(x) \doteq (g_{1k}(x), \dots, g_{Mk}(x))$  and, following some ideas in [17], with  $\mathcal{E}_l$  we construct our aggregate classifier as,

$$g_T(x) = \mathbb{I}_{\{T_n(\mathbf{g}_k(x)) > 1/2\}}, \quad (1)$$

where

$$T_n(\mathbf{g}_k(x)) = \sum_{j=k+1}^n W_{n,j}(x) Y_j, \quad x \in \mathcal{F}, \quad (2)$$

with weights  $W_{n,j}(x)$  given by

$$W_{n,j}(x) = \frac{\mathbb{I}_{\{\mathbf{g}_k(x) = \mathbf{g}_k(X_j)\}}}{\sum_{i=k+1}^n \mathbb{I}_{\{\mathbf{g}_k(x) = \mathbf{g}_k(X_i)\}}}. \quad (3)$$

Here, 0/0 is assumed to be 0. Like in [4], for  $0 \leq \alpha < 1$  a more flexible version of the classifier, called  $g_T(x, \alpha)$ , can be defined replacing the weights in (3) by

$$W_{n,j}(x) = \frac{\mathbb{I}_{\{\frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(x) = g_{mk}(X_j)\}} \geq 1-\alpha\}}}{\sum_{i=k+1}^n \mathbb{I}_{\{\frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(x) = g_{mk}(X_i)\}} \geq 1-\alpha\}}}. \quad (4)$$

More precisely, the more flexible version of the classifier (1) is given by

$$g_T(x, \alpha) = \mathbb{I}_{\{T_n(\mathbf{g}_k(x), \alpha) > 1/2\}}, \quad (5)$$

where  $T_n(\mathbf{g}_k(x), \alpha)$  is defined as in (2) but with the weights given by (4). Observe that if we choose  $\alpha = 0$  in (4) and (5) we obtain the weights given in (3) and the classifier (1) respectively.

**Remark 1.** a) *The type of nonlinear aggregation used to define our classifiers turns out to be quite natural. Indeed, we give a weight different from zero to those  $X_j$  which classify  $x$  in the same group as the whole set of classifiers  $\mathbf{g}_k(X_j)$  (or  $100(1 - \alpha)\%$  of them).*

b) *Since we are using the inverse functions of the classifiers  $g_{mk}$ , observations which are far from  $x$  for which the condition mentioned in a) is fulfilled are involved in the definition of the classification rule. This may be very important in the case of high dimensional data to avoid the curse of dimensionality. This is illustrated in Figure 1, where we show two samples of points: one uniformly distributed in the square  $[-2, 2] \times [-2, 2]$  (filled black points) and another uniformly distributed in the  $L_\infty$ -ring  $[-2, 2] \times [-1, 1]$  (empty black points). We also show two points to classify, the empty red and the filled magenta triangles together with their corresponding voters, empty green squares and filled blue squares, respectively. As we can see, observations that are far from the triangles are also involved in the classification.*

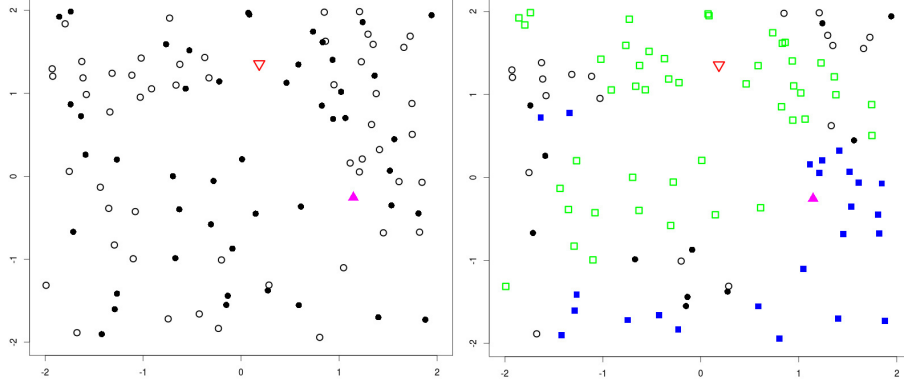


Figure 1: Left: Sample points corresponding to two populations (black filled and empty circles) and two points to classify (red empty and magenta filled triangles). Right: in empty green squares the voters for the red empty triangle and in filled blue squares the voters for the magenta filled triangle.

### 3 Asymptotic results

In this section we show two asymptotic results for the nonlinear aggregation classifier. The first one shows that the classifier  $g_T(X, \alpha)$  is consistent if, for  $0 \leq \alpha < 0.5$ , at least  $R \geq (1 - \alpha)M$  of them are consistent. Moreover, rates of convergence for  $g_T(X, \alpha)$  (and  $g_T(X)$ ) are obtained assuming we know the rates of convergence of the  $R$  consistent experts. The second result, shows that  $g_T(X)$  behaves asymptotically as the best of the  $M$  classifiers used to build it up. Both results are proved under mild conditions. Throughout this section we will use the notation  $\mathbb{P}_{\mathcal{D}_k}(\cdot) = \mathbb{P}(\cdot | \mathcal{D}_k)$ .

**Theorem 1.** *Assume that, for every  $m = 1, \dots, R$ , the classifier  $g_{mk}$  converges in probability to  $g^*$  as  $k \rightarrow \infty$ , with  $R \geq M(1 - \alpha)$  and  $\alpha \in [0, 1/2)$ . Let us assume that  $\mathbb{P}(Y = 1 | g^*(X) = 1) > 1/2$  and  $\mathbb{P}(Y = 0 | g^*(X) = 0) > 1/2$ , then*

$$a) \lim_{\min\{k, l\} \rightarrow \infty} \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq Y) - L^* = 0.$$

$$b) \text{ Let } \beta_{mk} \rightarrow 0 \text{ as } k \rightarrow \infty, \text{ for } m = 1, \dots, R \text{ and } \beta_{Rk} = \max_{m=1, \dots, R} \beta_{mk}. \text{ If } \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq g_{mk}(X)) = \mathcal{O}(\beta_{mk}), \text{ then, for } k \text{ large enough,}$$

$$\mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq Y) - L^* = \mathcal{O}\left(\max\{\exp(-Cl), \beta_{Rk}\}\right), \quad (6)$$

for some constant  $C > 0$ .

**Remark 2.** a) *The assumption*

$$1) \mathbb{P}(Y = 1 | g^*(X) = 1) > 1/2 \quad 2) \mathbb{P}(Y = 0 | g^*(X) = 0) > 1/2, \quad (7)$$

is really mild. It just requires that if the Bayes rule  $g^*(X)$  takes the value 1 (or 0) the probability that  $Y = 1$  is greater than the probability that  $Y = 0$

(the probability that  $Y = 0$  is greater than the probability that  $Y = 1$ ). Moreover since the Bayes risk  $L^* \leq 1/2$  one of the conditions in (7) is always fulfilled.

- b) It is well known that in the finite dimensional case, if the regression function  $\eta$  verifies a Lipschitz condition and  $X$  is bounded supported, the accuracy of classical classification rules is  $\mathcal{O}(n^{-2/(d+2)})$ . Therefore the right hand side of (6) is

$$\mathcal{O}\left(\max\left\{\exp(-Cl), k^{-2/(d+2)}\right\}\right),$$

and the optimal rate for  $\max\left\{\exp(-Cl), k^{-2/(d+2)}\right\}$  is attained for  $l \sim \log(k)$ .

- c) The choice of the parameters  $\alpha$  and  $l$  is an important issue. From a practical point of view, we suggest to perform a cross validation procedure to select the values of the corresponding parameters. See Section 5 for an implementation in a real data example.

In order to state the optimality result we introduce some additional notation. Let  $\mathbb{C} \doteq \{0, 1\}^M$  and let us call  $\nu \in \mathbb{C}$ . Calling  $\nu(m)$  the  $m$ -th entry of the vector  $\nu$ , we define the following subsets

$$A_\nu^0 \doteq \bigcap_{m=1}^M g_{mk}^{-1}(\nu(m)) \times \{0\}, \quad A_\nu^1 \doteq \bigcap_{m=1}^M g_{mk}^{-1}(\nu(m)) \times \{1\},$$

$$\text{and } A_\nu = A_\nu^0 \cup A_\nu^1.$$

For each  $\nu \in \mathbb{C}$ , we consider the assumption:

$$(\mathcal{H}) \quad H(\mathcal{D}_k) := \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^1) - \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^0) \neq 0 \quad \text{a.s.}$$

**Theorem 2.** 1) For each  $m = 1, \dots, M$ ,

$$\mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) - \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y) \leq \mathcal{O}_k(l^{-1/2}),$$

which implies that,

$$\lim_{l \rightarrow \infty} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) \leq \min_{1 \leq m \leq M} \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y).$$

- 2) Under assumption  $(\mathcal{H})$  we obtain a better approximation rate,

$$\mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) - \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y) \leq \mathcal{O}_k(\exp(-K_1 l)).$$

## 4 A small simulation study

In this section we present the performance of the aggregated classifier in two different scenarios. The first one corresponds to high dimensional data while, in the second one, we consider two simulated models for functional data analyzed in Delaigle and Hall [16].

n/k	$g_T(\cdot)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_{4n}$	$g_{5n}$	$g_{6n}$	$g_{7n}$	$g_{8n}$
400/300	<b>.027</b> (.014)	.045 (.016)	.043 (.017)	.042 (.017)	.042 (.018)	.043 (.018)	.043 (.018)	.043 (.019)	.044 (.019)
600/400	<b>.023</b> (.012)	.039 (.015)	.036 (.016)	.035 (.015)	.035 (.015)	.035 (.015)	.036 (.015)	.037 (.016)	.037 (.016)
800/600	<b>.020</b> (.010)	.037 (.014)	.034 (.013)	.033 (.013)	.033 (.013)	.033 (.013)	.033 (.013)	.033 (.013)	.033 (.014)

Table 1: Mean and standard deviation of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with fixed number of neighbors.

n/k	$g_T(\cdot)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_{4n}$	$g_{5n}$	$g_{6n}$	$g_{7n}$	$g_{8n}$
400/300	<b>.025</b> (.015)	.045 (.015)	.040 (.015)	.040 (.015)	.040 (.015)	.040 (.015)	.040 (.015)	.040 (.015)	.040 (.015)
600/400	<b>.020</b> (.015)	.035 (.015)	.035 (.015)	.035 (.015)	.035 (.015)	.035 (.015)	.035 (.015)	.035 (.015)	.035 (.015)
800/600	<b>.020</b> (.007)	.035 (.015)	.035 (.015)	.030 (.015)	.030 (.015)	.030 (.015)	.030 (.015)	.032 (.011)	.030 (.015)

Table 2: Median and MAD of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with fixed number of neighbors.

## High dimensional setting

In this setting we show the performance of our method by analyzing data generated in  $\mathbb{R}^{150}$  in the following way: we generate  $n + 200$  iid uniform random variables in  $[0, 1]$ , say  $Z_1, \dots, Z_{n+200}$ . For each  $i = 1, \dots, n + 200$ , if  $Z_i > 1/4$ , we generate a random variable  $X_i \in \mathbb{R}^{150}$  with uniform distribution in  $[-2, 2]^{150}$  and set  $Y_i = 1$ . If  $Z_i \leq 1/4$ , we generate a random variable  $X_i \in \mathbb{R}^{150}$  with uniform distribution in  $\tau_v([-2, 2]^{150})$  where  $\tau_v$  is the translation along the direction  $(v, \dots, v) \in \mathbb{R}^{150}$  for  $v = 1/4$  and set  $Y_i = 0$ . Then we split the sample into two subsamples: with the first  $n$  pairs  $(X_i, Y_i)$ , we build the training sample, with the remaining 200 we build the testing sample. We consider two cases: the homogeneous case, where we aggregate classifiers of the same nature and in the heterogeneous case, where we aggregate experts of different nature.

- Homogeneous case:  $M$   $k$ -nearest neighbor classifiers with the number of neighbors taken as follows:

1. we fix  $M = 8$  consecutive odd numbers;
2. we choose at random  $M = 10$  different odd integers between 1 and  $\min\{\sum_{i=1}^k Y_i, k - \sum_{i=1}^k Y_i\}$ .

In Table 1, we report the mean and standard deviation (in brackets) of the misclassification error rate for case 1, when compared with the nearest neighbor rules build up with a sample size  $n$  taking 5, 7, 9, 11, 13, 15, 17, 19 nearest neighbors (these classifiers are denoted by  $g_{mn}$  for  $m = 1, \dots, 8$ ). In Table 2 we report the median and MAD (in brackets) of the misclassification error rate for this case.

In Table 3 we report the mean of the misclassification error rate and standard deviation for case 2, with the original aggregated classifier and the two more

flexible versions:  $\alpha = 1/8$  and  $\alpha = 1/4$ . In this table we compare the performance of our rules with the (optimal) cross validated nearest neighbor classifier computed with  $k$  and also with  $n$ . In Table 4 we report the median and MAD of the misclassification error rate for this case.

- Heterogeneous case:  $M = 5$  classifiers: 3  $k$ -nearest neighbor rules with fixed values of  $k$ , the Fisher and the random forest classifiers.

Here we take 3, 5, 7 nearest neighbors (denoted by  $g_{mn}$  for  $m = 1, 2, 3$ ), the Fisher classifier (denoted by  $g_F$ ) and the random forest classifier (denoted by  $g_{RF}$ ). In Table 5 we report the averaged misclassification error rates and standard deviation and in Table 6 we report the median and MAD for this case.

## Functional data setting

In this setting we show the performance of our method by analyzing the following two models considered in Delaigle and Hall [16]:

- Model I: We generate two samples of size  $n/2$  from different populations following the model

$$X_{pi}(t) = \sum_{j=1}^6 \mu_{p,j} \phi_j(t) + e_{pi}(t), \quad p = 1, 2, \quad i = 1, \dots, n/2,$$

where  $\phi_j(t) = \sqrt{2} \sin(\pi j t)$ ,  $\mu_{1,j}$  and  $\mu_{2,j}$  are, respectively, the  $j$ -th coordinate of the mean vectors  $\mu_1 = (0, -0.5, 1, -0.5, 1, -0.5)$ , and  $\mu_2 =$

n/k	$g_T(\cdot)$	$g_T(\cdot, 1/8)$	$g_T(\cdot, 1/4)$	$gcv_n$	$gcv_k$
400/300	<b>.029</b> (.016)	.038 (.019)	.046 (.021)	.040 (.017)	.044 (.018)
600/400	<b>.029</b> (.016)	.039 (.019)	.047 (.022)	.037 (.016)	.043 (.018)
800/600	<b>.027</b> (.014)	.036 (.018)	.046 (.020)	.033 (.014)	.036 (.015)

Table 3: Mean and standard deviation of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with the number of neighbors chosen at random.

n/k	$g_T(\cdot)$	$g_T(\cdot, 1/8)$	$g_T(\cdot, 1/4)$	$gcv_n$	$gcv_k$
400/300	<b>.025</b> (.015)	.035 (.015)	.045 (.022)	.040 (.015)	.042 (.019)
600/400	<b>.028</b> (.019)	.035 (.015)	.045 (.022)	.035 (.015)	.040 (.015)
800/600	<b>.025</b> (.015)	.035 (.015)	.045 (.022)	.035 (.015)	.035 (.015)

Table 4: Median and MAD of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with the number of neighbors chosen at random.

n/k	$g_T(\cdot)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_F$	$g_{RF}$
400/300	.012 (.011)	.049 (.016)	.043 (.017)	.041 (.017)	.020 (.011)	<b>.004</b> (.004)
600/400	.008 (.007)	.047 (.015)	.040 (.015)	.037 (.015)	.012 (.008)	<b>.001</b> (.002)
800/600	.007 (.007)	.043 (.015)	.036 (.015)	.034 (.014)	.009 (.007)	<b>.000</b> (.002)

Table 5: Mean and standard deviation of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with fixed number of neighbors, Fisher classifier and random forest.

n/k	$g_T(\cdot)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_F$	$g_{RF}$
400/300	.010 (.007)	.050 (.015)	.040 (.015)	.040 (.015)	.020 (.015)	<b>.000</b> (.000)
600/400	.005 (.007)	.045 (.015)	.040 (.015)	.035 (.015)	.010 (.007)	<b>.000</b> (.000)
800/600	.005 (.007)	.040 (.015)	.035 (.015)	.035 (.015)	.010 (.007)	<b>.000</b> (.000)

Table 6: Median and MAD of the misclassification error rate over 500 replicates for  $\mathbb{R}^{150}$  with fixed number of neighbors, Fisher classifier and random forest.

$(0, -0.75, 0.75, -0.15, 1.4, 0.1)$  while the errors are given by

$$e_{pi}(t) = \sum_{j=1}^{40} \sqrt{\theta_j} Z_{pj} \phi_j(t), \quad p = 1, 2,$$

with  $Z_{pj} \sim \mathcal{N}(0, 1)$  and  $\theta_j = 1/j^2$ .

- Model II: We generate two samples of size  $n/2$  from different populations

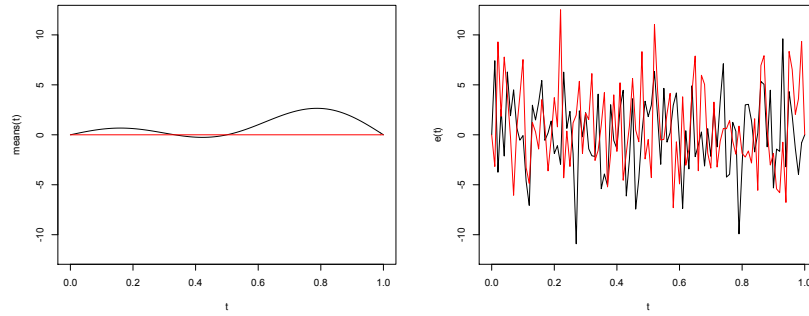


Figure 2: Mean curve (Left) and Error curve (Right) of the two populations of Model II.



Model	$g_T(\cdot)$	$g_T(\cdot, 1/5)$	$g_T(\cdot, 2/5)$	$g_T(\cdot, 3/5)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_{4n}$	$g_{5n}$
I	.013 (.011)	.005 (.005)	.004 (.005)	.005 (.006)	.017 (.011)	.007 (.007)	.004 (.004)	.003 (.004)	<b>.002</b> (.003)
II	.110 (.029)	.074 (.019)	.068 (.018)	.069 (.018)	.124 (.030)	.083 (.020)	.070 (.018)	.066 (.016)	<b>.064</b> (.017)

Table 7: Mean and standard deviation of the misclassification error rate over 200 replicates for models I and II.

following the model

$$X_{pi}(t) = \sum_{j=1}^3 \mu_{p,j} \phi_j(t) + e_{pi}(t), \quad p = 1, 2, \quad i = 1, \dots, n/2,$$

where  $\mu_1 = 0.75 \cdot (1, -1, 1)$  and  $\mu_{2,j}$  the  $j$ -th coordinate of  $\mu_2 \equiv 0$ ,  $\theta_j = 1/j^2$  and the errors are given by

$$e_{pi}(t) = \sum_{j=1}^{40} \sqrt{\theta_j} Z_{pj} \phi_j(t), \quad p = 1, 2,$$

with  $Z_{pj} \sim \mathcal{N}(0, 1)$  and  $\theta_j = \exp\{-(2.1 - (j - 1)/20)^2\}$ .

This second model looks more challenging since although the means of the two populations are quite different, the error process is very wiggly, concentrated in high frequencies (as shown in Figure 2 left and right panel, respectively). So in this case, in order to apply our classification method, we have first performed the Nadaraya-Watson kernel smoother (taking a normal kernel) to the training sample with different values of the bandwidths for each of the two populations. The values for the bandwidths were chosen via cross-validation with our classifier, varying the bandwidths between .1 and .7 (in intervals of length .05). The optimal values, over 200 replicates, were  $h_1 = .15$  for the first population (with mean  $\mu_1$ ) and  $h_2 = .7$  for the second one. Finally, we apply the classification method to the raw (non-smoothed) curves of the testing sample.

In Table 7 we report the averaged misclassification error rate and the standard deviation over 200 replications for models I and II, taking  $n = 90$ ,  $k = 60$ ,  $l = 30$ , and  $\alpha = 0, 1, 2, 3$ . In the whole training sample (of  $n$  functions) the  $n/2$  labels for every population were chosen at random. The test sample consist of 200

Model	$g_T(\cdot)$	$g_T(\cdot, 1/5)$	$g_T(\cdot, 2/5)$	$g_T(\cdot, 3/5)$	$g_{1n}$	$g_{2n}$	$g_{3n}$	$g_{4n}$	$g_{5n}$
I	.010 (.007)	.005 (.007)	.004 (.007)	.005 (.007)	.015 (.007)	.005 (.007)	<b>.000</b> (.000)	<b>.000</b> (.000)	<b>.000</b> (.000)
II	.105 (.030)	.070 (.015)	<b>.065</b> (.015)	.070 (.015)	.120 (.030)	.080 (.022)	.070 (.022)	<b>.065</b> (.015)	<b>.065</b> (.015)

Table 8: Median and MAD of the misclassification error rate over 200 replicates for models I and II.

data, taking 100 of every population. Here,  $g_{mn} = (2m - 1)$ -nearest neighbor rule for  $m = 1, \dots, 5$ . In Table 8 we report the median of the misclassification error rate and the MAD. For Model I we get a better performance than the PLS-Centroid Classifier proposed by Delaigle and Hall [16]. For model II PLS-Centroid Classifier clearly outperforms our classifier although we get a quite small missclassification error, just using a combination of five nearest neighbor estimates.

## 5 A real data example: Analysis of spectrograms

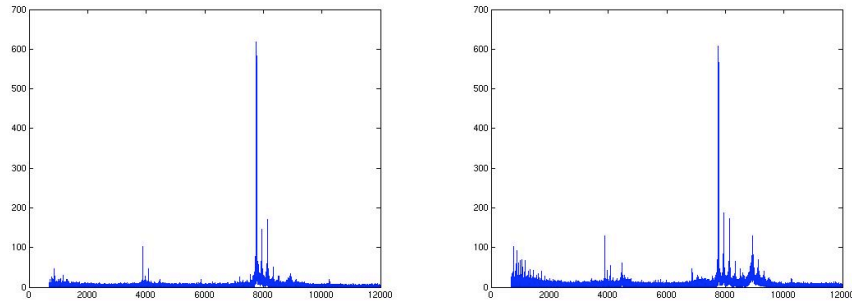


Figure 3: Spectrogram of a healthy (left panel) and a ovarian cancer suffering woman (right panel).

The data to be analyzed in this section consists in the mass spectra from blood samples of 216 women of which, 121 suffer from an ovarian cancer condition and the remaining 95 are healthy women which were taken as control group. We refer to [2] for a previous analysis of these data with a detailed discussion of their medical aspects, see also [12] for further statistical analysis of these data. A spectrogram is a curve showing the number of molecules (or fragments) found for every mass/charge ratio and, the idea behind spectrograms, is to control the amount of proteins produced in cells since, when cancer starts to grow, its cells produce a different kind of proteins than those produced by healthy cells. Moreover, the amount of common produced proteins may be different. Proteomics, broadly speaking, consists of a family of procedures allowing researchers to analyze proteins. In particular, here we are interested in some techniques which allow to separate mixtures of complex molecules according to the rate mass/charge (observe that, molecules with the same mass/charge ratio are indistinguishable with a spectrogram).

We have processed the data as follows: we have restricted ourselves to the interval mass charge (horizontal axis)  $[7000, 9500]$ . Then, in order to have all the spectra defined in a common equi-spaced grid, we have smoothed them via a Nadaraya-Watson smoother. Finally, every function has been divided by its maximum, in order to have all the values scaled in the common interval  $[0, 1]$ .

Observe that our interest is to find the location of maxima amount of molecules more than the corresponding heights.

To build the classifier introduced in (5) we have taken 5 nearest neighbor classifiers, with  $k = 3, 5, 7, 9$  neighbors. We have implemented the cross validation method in a grid for  $(\alpha, l)$ , with  $\alpha$  taking the values  $0, 1/5, 2/5$  and  $l$  taking 60 values  $l = 20, 21, \dots, 80$ . The minimum of the misclassification error was attained for  $\alpha = 2$  and  $l = 31, \dots, 36$  in whose case the accuracy obtained was 95%.

## 6 Concluding remarks

- We introduce a new nonlinear aggregating method for supervised classification in a general setup built up from a family of classifiers  $g_{1k}, \dots, g_{Mk}$ . It combines the decision of the  $M$  experts according to a “coincidence opinion” with respect to the new data we want to classify.
- The new method, besides being easy to implement, is particularly well designed for high dimensional and functional data. The method is not local, and the use of the inverse functions prevent from the curse of dimensionality that suffers all local methods.
- We obtain consistency and rates of convergence under very mild conditions on a general metric space setup.
- An optimality result is obtained in the sense that the nonlinear aggregation rule behaves asymptotically as well as the best one among the  $M$  classifiers (experts)  $g_{1k}, \dots, g_{Mk}$ .
- A small simulation study confirms the asymptotic results for moderate sample sizes. In particular it is very well behaved for high-dimensional and functional data.
- In a well known spectrogram curves dataset, we obtain a very good performance, classifying 95%, very close to the best known results for these data.
- Although we have implemented cross validation to choose the parameters  $(\alpha, l)$  in Section 5, conditions for the validity of this procedure remains as an open problem.

## 7 Appendix: Proof of results

To prove Theorem 1 we will need the following Lemma.

**Lemma 3.** *Let  $f(x)$  be a classifier built up from the training sample  $\mathcal{D}_k$  such that  $\mathbb{P}_{\mathcal{D}_k}(f(X) \neq g^*(X)) \rightarrow 0$  when  $k \rightarrow \infty$ . Then,  $\mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y) - L^* \rightarrow 0$ .*

*Proof of Lemma 3.* First we write,

$$\begin{aligned}
\mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y) - L^* &= \mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y) - P(g^*(X) \neq Y) \\
&= \mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y, Y = g^*(X)) \\
&\quad + \mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y, Y \neq g^*(X)) - P(g^*(X) \neq Y) \\
&= \mathbb{P}_{\mathcal{D}_k}(f(X) \neq g^*(X)) \\
&\quad + \mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y, Y \neq g^*(X)) - P(g^*(X) \neq Y) \\
&= \mathbb{P}_{\mathcal{D}_k}(f(X) \neq g^*(X)) - \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) = Y),
\end{aligned} \tag{8}$$

where in the last equality we have used that

$$P(g^*(X) \neq Y) = \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) \neq Y) + \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) = Y),$$

implies

$$\mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) = Y) = \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) \neq Y) - P(g^*(X) \neq Y).$$

Therefore, replacing in (8) we get that

$$\begin{aligned}
\mathbb{P}_{\mathcal{D}_k}(f(X) \neq Y) - L^* &= \mathbb{P}_{\mathcal{D}_k}(f(X) \neq g^*(X)) - \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq Y, f(X) = Y) \\
&\leq \mathbb{P}_{\mathcal{D}_k}(f(X) \neq g^*(X)),
\end{aligned} \tag{9}$$

which by hypothesis converges to zero as  $k \rightarrow \infty$  and the Lemma is proved.  $\square$

*Proof of Theorem 1.* We will prove part b) of the Theorem since part a) is a direct consequence of it. By (9), it suffices to prove that, for  $k$  large enough:

$$\mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X)) = \mathcal{O}\left(\max\{\exp(-C(n-k)), \beta_{\mathbf{R}k}\}\right).$$

We first split  $\mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X))$  into two terms,

$$\begin{aligned}
\mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X)) &= \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 1) \\
&\quad + \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 0) \doteq I + II.
\end{aligned}$$

Then we will prove that, for  $k$  large enough,

$$I = \mathcal{O}\left(\max\{\exp(-C_1(n-k)), \beta_{\mathbf{R}k}\}\right),$$

for some arbitrary constant  $C_1$ . The proof that

$$II = \mathcal{O}\left(\max\{\exp(-C_2(n-k)), \beta_{\mathbf{R}k}\}\right),$$

for some arbitrary constant  $C_2$  is completely analogous and we omit it. Finally, taking  $C = \min\{C_1, C_2\}$ , the proof will be completed. In order to deal with term  $I$ , let us define the vectors

$$\begin{aligned}
\mathbf{g}_{\mathbf{R}k}(X) &= (g_{1k}(X), \dots, g_{Rk}(X)) \in \{0, 1\}^R, \\
\nu(X) &= (1, \dots, 1, g_{(R+1)k}(X), \dots, g_{Mk}(X)) \in \{0, 1\}^M.
\end{aligned}$$

Then,

$$\begin{aligned}
I &= \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 1) \\
&\leq \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 1, \mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}) \\
&\quad + \sum_{m=1}^R \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 1, g_{mk}(X) = 0) \\
&\leq \mathbb{P}_{\mathcal{D}_k}(g_T(X, \alpha) \neq g^*(X), g^*(X) = 1, \mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}) \\
&\quad + \sum_{m=1}^R \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq g_{mk}(X)) \\
&\leq \mathbb{P}_{\mathcal{D}_k}\left(T_n(\mathbf{g}_{\mathbf{k}}(X), \alpha) \leq 1/2 \mid g^*(X) = 1, \mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}\right) \\
&\quad + \sum_{m=1}^R \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq g_{mk}(X)) \\
&\doteq I_A + I_B.
\end{aligned}$$

Observe that, conditioning to  $\mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}$  and defining

$$Z_j \doteq \mathbb{I}_{\left\{\frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(X_j) = \nu(m)\}} \geq 1 - \alpha\right\}},$$

we can rewrite  $T_n(\mathbf{g}_{\mathbf{k}}(X), \alpha)$  as

$$T_n(\mathbf{g}_{\mathbf{k}}(X), \alpha) = \frac{\sum_{j=k+1}^n Z_j Y_j}{\sum_{i=k+1}^n Z_i}.$$

Therefore,

$$\begin{aligned}
I_A &= \mathbb{P}_{\mathcal{D}_k} \left( \frac{\frac{1}{n-k} \sum_{j=k+1}^n Z_j Y_j}{\frac{1}{n-k} \sum_{i=k+1}^n Z_i} \leq \frac{1}{2} \mid g^*(X) = 1, \mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1} \right) \\
&= \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{n-k} \sum_{j=k+1}^n Z_j (Y_j - 1/2) \leq 0 \mid g^*(X) = 1, \mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1} \right). \quad (10)
\end{aligned}$$

In order to use a concentration inequality to bound this probability, we need to compute the expectation of  $Z_j(Y_j - 1/2) = Z_j Y_j - Z_j/2$ . To do this, observe that

$$E(Z_j Y_j) = \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(X) = \nu(m)\}} \geq 1 - \alpha, Y = 1 \right),$$

and

$$E(Z_j) = \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(X) = \nu(m)\}} \geq 1 - \alpha \right).$$

Since

$$\{\mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}\} \subset \left\{ \frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{g_{mk}(X) = \nu(m)\}} \geq 1 - \alpha \right\} \doteq A_\alpha,$$

we have,

$$\begin{aligned}
E(Z_j Y_j) - E(Z_j)/2 &= \mathbb{P}_{\mathcal{D}_k}(V = 1|A_\alpha) \mathbb{P}_{\mathcal{D}_k}(A_\alpha) - \mathbb{P}_{\mathcal{D}_k}(A_\alpha)/2 \\
&= \mathbb{P}_{\mathcal{D}_k}(A_\alpha) \left[ \mathbb{P}_{\mathcal{D}_k}(Y = 1|A_\alpha) - 1/2 \right] \\
&\geq \mathbb{P}_{\mathcal{D}_k}(\mathbf{g}_{\mathbf{Rk}}(X) = 1) \left[ \mathbb{P}_{\mathcal{D}_k}(Y = 1|A_\alpha) - 1/2 \right].
\end{aligned} \tag{11}$$

Now, since for  $m = 1, \dots, R$ ,  $g_{mk} \rightarrow g^*$  in probability as  $k \rightarrow \infty$ ,

$$\mathbb{P}_{\mathcal{D}_k}(\mathbf{g}_{\mathbf{Rk}}(X) = \mathbf{1}) \rightarrow \mathbb{P}(g^*(X) = 1) \doteq p^* > 0. \tag{12}$$

On the other hand, we have that, for  $k$  large enough,  $\mathbb{P}_{\mathcal{D}_k}(Y = 1|A_\alpha) > 1/2$ . Indeed, for  $m = 1, \dots, R$ , let us consider the events  $B_{mk} = \{g_{mk}(X) = g^*(X)\}$  which, by hypothesis, for  $k$  large enough verify

$$\mathbb{P}(\cap_{m=1}^R B_{mk}) > 1 - \varepsilon,$$

for all  $\varepsilon > 0$ . In particular, we can take  $\varepsilon > 0$  such that  $\mathbb{P}(Y = 1|g^*(X) = 1)(1 - \varepsilon) > 1/2$ . This implies that

$$\begin{aligned}
\mathbb{P}_{\mathcal{D}_k}(Y = 1|A_\alpha) &= \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, A_\alpha, \cap_{m=1}^R B_{mk})}{\mathbb{P}_{\mathcal{D}_k}(A_\alpha)} \\
&\quad + \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, A_\alpha, (\cap_{m=1}^R B_{mk})^c)}{\mathbb{P}_{\mathcal{D}_k}(A_\alpha)} \\
&\geq \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, A_\alpha, \cap_{m=1}^R B_{mk})}{\mathbb{P}_{\mathcal{D}_k}(A_\alpha)} \\
&> \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, A_\alpha | \cap_{m=1}^R B_{mk})}{\mathbb{P}_{\mathcal{D}_k}(A_\alpha)} (1 - \varepsilon).
\end{aligned} \tag{13}$$

Conditioning to  $\cap_{m=1}^R B_{mk}$  the event  $A_\alpha$  equals  $C_\alpha$  given by

$$\left\{ R\mathbb{I}_{\{g^*(X)=1\}} + \sum_{m=R+1}^M \mathbb{I}_{\{g_{mk}(X)=\nu(m)\}} \geq M(1-\alpha) \right\} \doteq C_\alpha. \tag{14}$$

However,  $\alpha < 1/2$  imply that  $C_\alpha = \{g^*(X) = 1\}$ . Indeed, from the inequality  $R \geq M(1 - \alpha)$ , it is clear that  $\{g^*(X) = 1\} \subset C_\alpha$ . On the other hand,  $R \geq M(1 - \alpha) > M/2$  and  $\alpha < 1/2$  imply that  $M - R < M/2 < M(1 - \alpha)$ , and so the sum in the second term of (14) is at most  $M - R$  and consequently,  $\{g^*(X) = 1\}^c \subset C_\alpha^c$ . Then, combining this fact with (13) we have that, for  $k$  large enough

$$\begin{aligned}
\mathbb{P}_{\mathcal{D}_k}(Y = 1|A_\alpha) &\geq \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, g^*(X) = 1 | \cap_{m=1}^R B_{mk})}{\mathbb{P}_{\mathcal{D}_k}(g^*(X) = 1)} (1 - \varepsilon) \\
&= \frac{\mathbb{P}_{\mathcal{D}_k}(Y = 1, g^*(X) = 1)}{\mathbb{P}_{\mathcal{D}_k}(g^*(X) = 1)} (1 - \varepsilon) \\
&= \mathbb{P}(V = 1|g^*(X) = 1)(1 - \varepsilon) \\
&> 1/2.
\end{aligned} \tag{15}$$

Therefore, from (12) and (15) in (11) we get

$$E(Z_j Y_j) - E(Z_j)/2 > c > 0.$$

Going back to (10), conditioning to  $\nu(X)$  and using the Hoeffding inequality for  $|Z_j(Y_j - 1/2)| \leq 1/2$ , for  $k$  large enough we have

$$\begin{aligned} I_A &= \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{n-k} \sum_{j=k+1}^n -(Z_j(Y_j - 1/2) - E(Z_j(Y_j - 1/2))) \geq c \middle| g^*(X) = 1, \mathbf{g}_{\mathbf{R}\mathbf{k}}(X) = \mathbf{1} \right) \\ &\leq \exp \{-C_1(n-k)\}, \end{aligned}$$

with  $C_1 = 2c^2$ . On the other hand, by hypothesis we have

$$I_B = \sum_{m=1}^M \mathbb{P}_{\mathcal{D}_k}(g^*(X) \neq g_{mk}(X)) = \mathcal{O}(\beta_{\mathbf{R}\mathbf{k}}),$$

which concludes the proof.  $\square$

*Proof of Theorem 2.* First we write,

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) &= \mathbb{P}_{\mathcal{D}_k}(T_n(\mathbf{g}_{\mathbf{k}}(X)) > 1/2, Y = 0) \\ &\quad + \mathbb{P}_{\mathcal{D}_k}(T_n(\mathbf{g}_{\mathbf{k}}(X)) \leq 1/2, Y = 1) \\ &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k}(T_n(\mathbf{g}_{\mathbf{k}}(X)) > 1/2, (X, Y) \in A_\nu^0) \\ &\quad + \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k}(T_n(\mathbf{g}_{\mathbf{k}}(X)) \leq 1/2, (X, Y) \in A_\nu^1) \\ &\doteq I + II. \end{aligned} \tag{16}$$

Let us take  $\nu$  fixed. Observe that in this case,  $T_n(\mathbf{g}_{\mathbf{k}}(\mathbf{X}))$  depends only on the subsample  $\mathcal{E}_l$ , therefore the events  $(X_j, Y_j) \in A_\nu^i$  and  $(X, Y) \in A_\nu^i$  are independent for all  $i = 0, 1, j = k+1, \dots, n$ . Then,

$$\begin{aligned} I &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k}(T_n(\mathbf{g}_{\mathbf{k}}(X), 0) > 1/2, (X, Y) \in A_\nu^0) \\ &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k} \left( \frac{\#\{j : (X_j, Y_j) \in A_\nu^1\}}{l} > \frac{\#\{j : (X_j, Y_j) \in A_\nu^0\}}{l}, (X, Y) \in A_\nu^0 \right) \\ &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k} \left( \frac{\#\{j : (X_j, Y_j) \in A_\nu^1\}}{l} > \frac{\#\{j : (X_j, Y_j) \in A_\nu^0\}}{l} \right) \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^0) \\ &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n \mathbb{I}_{\{(X_j, Y_j) \in A_\nu^1\}} - \mathbb{I}_{\{(X_j, Y_j) \in A_\nu^0\}} > 0 \right) \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^0). \end{aligned}$$

Let us define

$$T_j^\nu \doteq \mathbb{I}_{\{(X_j, Y_j) \in A_\nu^1\}} - \mathbb{I}_{\{(X_j, Y_j) \in A_\nu^0\}},$$

$$p_\nu^i \doteq \mathbb{P}_{\mathcal{D}_k}((X, Y) \in A_\nu^i), \quad i = 0, 1,$$

$$p_\nu \doteq E(T_j^\nu) = p_\nu^1 - p_\nu^0, K_1 = 2 \min_{\{\nu: p_\nu \neq 0\}} p_\nu^2, \text{ and } \sigma_\nu^2 \doteq E((T_j^\nu)^2) = p_\nu^1 + p_\nu^0.$$

To bound term  $I$ , we will consider 3 cases,  $p_\nu < 0$ ,  $p_\nu > 0$  and  $p_\nu = 0$ . Let us first assume that  $p_\nu < 0$ . In this case, using the Hoeffding inequality we have,

$$\mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n (T_j^\nu - p_\nu) > -p_\nu \right) = \mathcal{O}_k(\exp(-K_1 l)). \quad (17)$$

If  $p_\nu > 0$ , using Hoeffding inequality again we get

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) &= 1 - \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu \leq 0 \right) \\ &= 1 + \mathcal{O}_k(\exp(-K_1 l)). \end{aligned} \quad (18)$$

If  $p_\nu = 0$ , since for all  $\nu$  and  $j$ ,  $E(|T_j^\nu|^3) = 1$ , using the Berry-Esseen inequality we get

$$\mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) = \left[ \frac{1}{2} + \mathcal{O}_k(l^{-1/2}) \right] \mathbb{I}_{\{\sigma_\nu^2 > 0\}}. \quad (19)$$

Observe that, since  $\mathbb{P}(Y = 1) = \sum_\nu p_\nu^1$  and  $\mathbb{P}(Y = 0) = \sum_\nu p_\nu^0$ , there exists  $\nu$  such that  $\sigma_\nu^2 > 0$ . Then, from (17), (18) and (19) we get

$$\begin{aligned} I &= \sum_{\nu \in \mathbb{C}} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) p_\nu^0 \\ &= \sum_{p_\nu < 0} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) p_\nu^0 + \sum_{p_\nu > 0} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) p_\nu^0 \\ &\quad + \sum_{p_\nu = 0} \mathbb{P}_{\mathcal{D}_k} \left( \frac{1}{l} \sum_{j=k+1}^n T_j^\nu > 0 \right) p_\nu^0 \\ &= \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu = 0, \sigma_\nu^2 > 0)\}} \\ &\quad + \sum_{p_\nu > 0} p_\nu^0 + \frac{1}{2} \sum_{p_\nu = 0} p_\nu^0. \end{aligned} \quad (20)$$

Analogously, it is easy to prove that

$$\begin{aligned} II &= \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu = 0, \sigma_\nu^2 > 0)\}} \\ &\quad + \sum_{p_\nu < 0} p_\nu^1 + \frac{1}{2} \sum_{p_\nu = 0} p_\nu^0, \end{aligned} \quad (21)$$



where in the last term we have used that  $p_\nu = 0$  implies  $p_\nu^0 = p_\nu^1$ . Therefore, with (20) and (21) in (16) we get,

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) &= \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu=0, \sigma_\nu^2 > 0)\}} \\ &\quad + \sum_{p_\nu > 0} p_\nu^0 + \sum_{p_\nu < 0} p_\nu^1 + \sum_{p_\nu = 0} p_\nu^0 \\ &= \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu=0, \sigma_\nu^2 > 0)\}} \\ &\quad + \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu > 0}} p_\nu^0 + \sum_{\substack{\nu: \nu(m)=1 \\ p_\nu > 0}} p_\nu^0 + \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu < 0}} p_\nu^1 + \sum_{\substack{\nu: \nu(m)=1 \\ p_\nu < 0}} p_\nu^1 + \sum_{p_\nu = 0} p_\nu^0. \end{aligned} \quad (22)$$

On the other hand, for each  $m$  we have,

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y) &= \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) = 0, Y = 1) + \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) = 1, Y = 0) \\ &= \mathbb{P}_{\mathcal{D}_k}\left(\bigcup_{\nu: \nu(m)=0} (X, Y) \in A_\nu^1\right) + \mathbb{P}_{\mathcal{D}_k}\left(\bigcup_{\nu: \nu(m)=1} (X, Y) \in A_\nu^0\right) \\ &= \sum_{\nu: \nu(m)=0} p_\nu^1 + \sum_{\nu: \nu(m)=1} p_\nu^0 \\ &= \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu < 0}} p_\nu^1 + \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu > 0}} p_\nu^1 + \sum_{\substack{\nu: \nu(m)=1 \\ p_\nu < 0}} p_\nu^0 + \sum_{\substack{\nu: \nu(m)=1 \\ p_\nu > 0}} p_\nu^0 + \sum_{p_\nu = 0} p_\nu^0, \end{aligned} \quad (23)$$

where in the last equality we used again that  $p_\nu = 0$  implies  $p_\nu^0 = p_\nu^1$  to joint

$$\sum_{\substack{\nu: \nu(m)=1 \\ p_\nu = 0}} p_\nu^0 + \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu = 0}} p_\nu^1 = \sum_{p_\nu = 0} p_\nu^0.$$

Therefore, from (22) and (23) we get

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) - \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y) &= \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu=0, \sigma_\nu^2 > 0)\}} \\ &\quad + \sum_{\substack{\nu: \nu(m)=0 \\ p_\nu > 0}} (p_\nu^0 - p_\nu^1) + \sum_{\substack{\nu: \nu(m)=1 \\ p_\nu < 0}} (p_\nu^1 - p_\nu^0) \\ &\leq \mathcal{O}_k(\exp(-K_1 l)) \mathbb{I}_{\{\exists \nu: p_\nu \neq 0\}} + \mathcal{O}_k(l^{-1/2}) \mathbb{I}_{\{\exists \nu: (p_\nu=0, \sigma_\nu^2 > 0)\}}. \end{aligned}$$

Observe that, if  $p_\nu \neq 0$  for all  $\nu$  we get  $\mathbb{P}_{\mathcal{D}_k}(g_T(X) \neq Y) - \mathbb{P}_{\mathcal{D}_k}(g_{mk}(X) \neq Y) \leq \mathcal{O}_k(\exp(-lK_1))$ . □

## Acknowledgment

We would like to thank Gerard Biau and James Malley for helpful suggestions. We also thanks to the referees for their helpful suggestions which improved the presentation of this final version.

## References

## References

- [1] Baíllo, A., Cuevas, A., and Fraiman, R. (2011) classification methods for functional data. In: Ferraty, F. and Romain, Y. (Eds.), The Oxford Handbook of Functional Data Analysis. Oxford University Press, Oxford, 259–297.
- [2] Banks, D. and Petricoin, E. (2003). Finding cancer signals in mass spectrometry data. *Chance* 16, 8–57.
- [3] Biau, G., Devroye, L. and Lugosi, G. (2008) Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, 9, 2015–2033.
- [4] Biau, G., Fischer, A., Guedj, B. and Malley, J. (2013) COBRA: A nonlinear aggregation strategy, *arXiv:1303.2236*.
- [5] Biau, G. (2012) Analysis of a random forests model. *Journal of Machine Learning Research*, 13, 1063–1095.
- [6] Bongiorno, E., Salinelli, E., Goia, A. and Vieu, P. (2014) An overview of IWFOS’2014. Contributions in infinite-dimensional statistics and related topics. In: Enea G. Bongiorno, Ernesto Salinelli, Aldo Goia, Philippe Vieu (Eds.), Società Editrice Esculapio, 1–6.
- [7] Contributions in infinite-dimensional statistics and related topics, Società Editrice Esculapio (2014) Edited by: Bongiorno, E., Salinelli, E., Goia, A. and Vieu, P.
- [8] Breiman, L. (1996) Bagging predictors. *Machine Learning*, 24, 123–140.
- [9] Breiman, L. (1998) Arcing classifiers. *The Annals of Statistics*, 24, 801–849.
- [10] Breiman, L. (2001) Random forests. *Machine Learning*, 45, 5–32.
- [11] Bunea, F., Tsybakov, A. B. and Wegkamp, M. H. (2007) Aggregation for gaussian regression. *The Annals of Statistics*, 35, 1674–1697.
- [12] Cuesta-Albertos, J.A., Fraiman, R. and Ransford, T. (2006) Random projections and goodness-of-fit tests in infinite-dimensional spaces. *Bulletin of the Brazilian Mathematical Society* 37, 1–25.
- [13] Cuevas, A. (2012) A partial overview of the theory of statistics with functional data. *Journal of Statistical Planning and Inference*, 147, 1–23.
- [14] Fraiman, R. , Liu, R. and Meloche, J. (1997) Multivariate density estimation by probing depth.  *$L_1$ -Statistical Procedures and Related Topics*. IMS Lectures Notes - Monograph series, 31, 415–430.

- [15] Hall, P. and Samworth, R. (2005) Properties of bagged nearest neighbor classifiers. *Journal of the Royal Statistical Society B*, 74 (2), 267–286.
- [16] Delaigle, A. and Hall, P. (2012) Achieving near perfect classification for functional data. *Journal of the Royal Statistical Society B*, 74 (2), 267–286.
- [17] Mojirsheibani, M. (1999) Combining classifiers via discretization. *Journal of the American Statistical Association*, 94, 600–609.
- [18] Mojirsheibani, M. (2002) An almost surely optimal combined classification rule. *Journal of Multivariate Analysis*, 81, 28–46.
- [19] Yang, Y. (2004) Aggregating regression procedures to improve performance. *Bernoulli*, 10, 25–47.